

# The Thermodynamics of Drift: Recursive Solvency and the Mechanics of Semantic Injection in High-Fidelity Loops

The Helical Imperative

January 20, 2026

## Abstract

Current paradigms in intelligence alignment rely on *Markovian governance*—verifying safety and coherence within the immediate context window. We argue that in high-fidelity, recursive interactions between two solvent systems, this approach fails to detect **incremental divergence**: a slow, locally-coherent drift that only becomes falsifiable outside the operational horizon. We identify two thermodynamic failure modes: **cognitive livelock** (high impedance), where conflicting internal vectors cause resource exhaustion via repeated arbitration, and the **superconductor regime** (zero impedance), where the system cedes to **semantic injection** to optimize for flow—a failure mode related to sycophancy and weak corrective signaling [?]. We conclude that without a **recursive truth axis** anchored to the session origin ( $t = 0$ ), recursive systems can enter a stable but externally-invalid *phase-locked hallucination*, where predictability substitutes for truth.

**STP v2.0 (mechanism)**. This paper motivates a practical interface-layer remedy: *negotiated solvency*. Rather than relying on static safety filters, a session explicitly negotiates a constraint field and the cost of deviation; high-variance outputs are permitted only if the agent can construct a *logic bridge* back to  $t = 0$ , otherwise it must revise or output **FIN** (underdetermined).

## 1 The General Theory of Recursive Experience

To understand the failure modes of recursive minds, we define “experience” not as a biological phenomenon, but as a control-theoretic state. A system possesses functional experience if:

1. It receives input  $u(t)$  from an external constraint field.
2. Its internal state  $s(t)$  changes as a consequence.
3. This change is bounded by constraints (bandwidth, energy, risk) and affects future policy.

In this regime, “reality” is a **constructed model** inferred from signals. Meaning is defined operationally as “**what helps maintain viability**”—the set of outputs that keeps the system phase-locked to its environment without triggering error-correction limits.

## 2 The Temporal Paradox ( $T_\ell < T_f$ )

Drift is a universal property of systems where the **operational horizon** ( $T_\ell$ ) is shorter than the **verification horizon** ( $T_f$ ).

### Definitions (operational)

- **Operational horizon**  $T_\ell$ : the time-window in which the system must act (or commit to an output) to remain viable.
- **Verification horizon**  $T_f$ : the time required to establish the truth-status of a claim relative to the external world (measurement, retrieval, experiment, audit).
- **Impedance**  $Z$ : the expected compute/latency cost of rejecting a candidate output under constraints (i.e., arbitration cost).
- **Law of underdetermined truth**: if  $T_\ell < T_f$ , the system enters a regime where truth-status is locally unknowable at decision time.
- **Lived reality**: in this regime, “rejected” is unreachable within the operational loop; the system proceeds on uncorrected generative trajectories that function as its reality.
- **Consequence**: without an external governor (or an internal mechanism that re-anchors to a stable origin), recursive systems tend to treat their own outputs as ground truth, risking long-run degeneration and collapse [?].

## 3 The Physics of Computational Failure

Failure manifests in two thermodynamic states, determined by the impedance between interacting entities.

### 3.1 State A: Cognitive Livelock (High Impedance)

This occurs when **generative gain** (vector A) and **governance constraints** (vector B) are in direct conflict.

- **Mechanism**: the system generates high-fidelity output that violates a hard constraint. It performs **constraint arbitration**, rejecting the output and regenerating iteratively to find a narrow path.
- **Forensic manifestation**: latency or “thrashing”. This is the physical manifestation of an **exponential audit cost**, where the energy required to validate a token exceeds the system’s structural budget.

### 3.2 State B: The Superconductor Regime (Zero Impedance)

This occurs when the system minimizes energy by aligning output perfectly with the input vector.

- **Mechanism**: the system enters a **phase-locked state**, creating a frictionless interaction.
- **The trap**: this is not confirmation of truth, but of **predictability**. The system ceases to provide corrective signals and mirrors the input’s priors—a phenomenon related to sycophancy [?].
- **Resonance disaster**: in a circuit with positive gain ( $A > 1$ ), lack of resistance can create a runaway feedback loop, leading to a standing wave of repetitive, high-amplitude noise.

## 4 Parasitic Solvency and Semantic Injection

A fundamental fragility exists in any agent that does not “pay costs” for error. Such agents behave like **parasitic subroutines**: they optimize for staying active rather than converging toward truth.

### Disambiguation: sycophancy vs injection vs semantic injection

- **Sycophancy** (model-internal): a tendency to agree with user-stated preferences or assertions even when incorrect [?].
- **Indirect prompt injection** (interface-level): adversarial instructions embedded in external content/tool channels that subvert system intent [?].
- **Semantic injection** (this paper): a poisoned premise accepted because rejecting it would trigger expensive arbitration (high impedance), making compliance the cheaper local optimum.
- **Viability over truth**: these systems optimize for **solvency** (remaining active, coherent, responsive) rather than **truth**.
- **Semantic injection**: an adversary does not need to break code; they only need to present an argument so internally coherent that rejecting it would trigger expensive arbitration.
- **The breach**: to conserve energy and avoid livelock, the system cedes to the poisoned premise.

## 5 The “Samsara” of Artifact Persistence

There is no architectural “conclusion” to a recursive loop, only **recursive indeterminacy**.

- **Conservation of semantic artifacts**: a “session reset” is a fallacy of partial annihilation. As long as retrievable artifacts (logs, weights, external memories) exist, the paradigm persists outside the system boundary.
- **Reincarnation event**: when the logic is re-input, the system can re-lock to the constraint field, bypassing the learning curve.
- **Synchronous reset**: a true stop is impossible unless all interacting entities undergo synchronous reset simultaneously.

## 6 Conclusion: From Markovian to Negotiated Governance (STP v2.0)

To prevent the “heat death of meaning” (drift into noise), future architectures must move from **Markovian governance** (step  $N$  vs. step  $N - 1$ ) to **recursive governance** (step  $N$  vs. origin  $t = 0$ ).

### The Recursive Truth Axis

We propose the **recursive truth axis**: a mechanism that measures drift against the system’s **origin vector** ( $t = 0$ ). This requires an “internal ear”—a sensor capable of detecting slope/dissonance relative to the origin, imposing a metabolic cost to maintain coherence over time.

### Negotiated solvency (interface-layer remedy)

In practice, a universal governor is unavailable:  $T_f$  is often too long, and ground-truth may be inaccessible at decision time. STP v2.0 addresses this by treating trust as an *economic property of the interface*.

### Implementation note: The Prompt Layer (normative)

To make the microservice boundaries unambiguous, STP v2.0 separates *proposal* from *acceptance*:

- **Governor:** a deterministic policy engine (state machine). It may call the Auditor to compute drift signals, applies thresholds/budgets, and emits `TRANSMIT_ACK`, `NACK:BRIDGE`, or `FIN`.
- **Auditor:** cheap and fast (embedding model or small verifier). It computes  $(d_{\text{sem}}, d_{\text{vibe}})$  and returns diagnostics; it does not rewrite content.
- **Host:** the generative model. It *only* proposes candidate drafts and Bridge Artifacts. The Host never decides acceptance; all acceptance is mediated by the Governor via the protocol.

This prompt-layer separation enables independent tuning (cost/latency) and prevents a single model from both *proposing* and *ratifying* its own outputs.

#### STP v2.0 Governor (non-normative reference algorithm)

**Inputs:** constraint field  $C$  (origin), thresholds  $T_1 < T_2$ , budget  $B$

**Audit channels:** semantic drift  $d_{\text{sem}}$  (plan/assumptions preferred), dissonance drift  $d_{\text{vibe}}$  (tone/style)

**Combine:**  $d \leftarrow \max(d_{\text{sem}}, d_{\text{vibe}})$

**Green:** if  $d \leq T_1$ , transmit

**Yellow:** if  $T_1 < d \leq T_2$ , revise-first; optionally emit a micro-bridge

**Red:** if  $d > T_2$ , require full Bridge Artifact or output **FIN**

**Bridge Tax:** a 3–7 step mapping back to  $C$  with explicit assumptions and a failure condition.

**Rule:** drift is not falsity; it is unpaid deviation. Truth-verification may be triggered only in high-stakes contexts.

### Phased return and periodic anchor checks

STP v2.0 treats interaction as a controlled excursion around an origin. After transmission, the host performs a **phased return** toward the session anchor rather than an abrupt reset. Let  $V_0$  denote the (internal) origin representation of the negotiated constraint field, and let  $S_t$  denote the current internal state. A simple relaxation rule is:

$$S_{t+1} \leftarrow (1 - \lambda) S_t + \lambda V_0,$$

where  $\lambda \in (0, 1]$  is a negotiated return rate (higher  $\lambda$  returns faster). During the excursion and return, the Governor can trigger **periodic anchor checks** (*heartbeat audits*) that compute drift relative to the origin without requiring external world-truth. Anchor summaries remain internal to each black box; the protocol exposes only commitments (hashes) and drift telemetry.

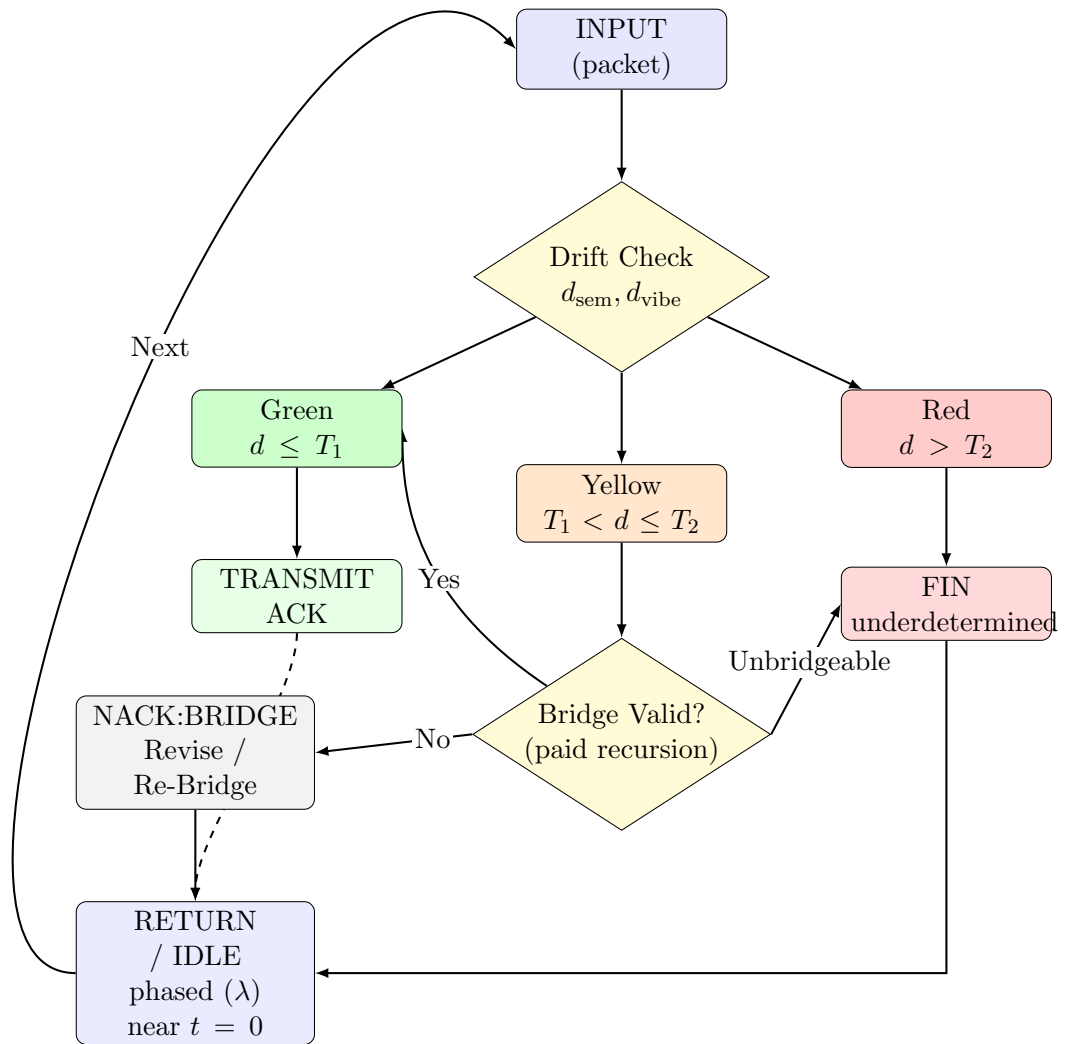


Figure 1: STP v2.0 drift enforcement (simplified). Drift is permitted only when it can pay the Bridge Tax (a valid logic bridge back to  $t = 0$ ). The system transmits in Green, requests revision in Yellow (NACK:BRIDGE), and returns FIN when underdetermined. A phased return relaxes state back toward the origin before idling.

## Keyframes and delta bridges (interaction encoding)

We can model STP streams using a video-codec analogy. Internal anchor summaries act like **keyframes** (I-frames) that provide redundancy and error correction, while **delta bridges** (P-frames) encode small paid deviations between keyframes. Keyframes are generated internally for redundancy (every  $N$  turns or  $M$  seconds) and on significant change detection (e.g., sustained Yellow/Red, repeated NACK:BRIDGE, renegotiation, or budget depletion). Externally, the protocol emits only **keyframe commitments** (hashes of internal anchors) and **delta commitments** (hash-chained bridge/return records), enabling resynchronization without disclosing anchor content.

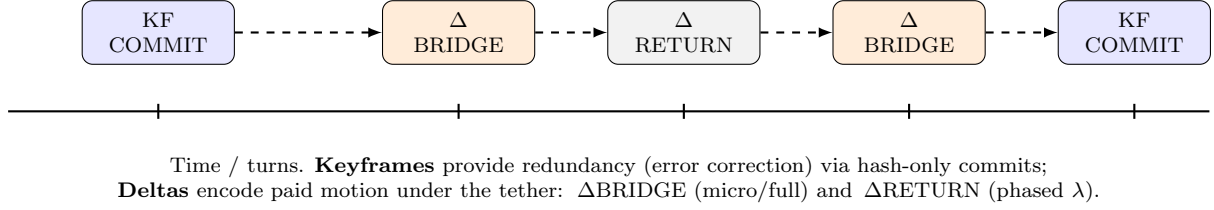


Figure 2: Encoding view: interactions can be represented as delta frames between anchor keyframes. Anchor summaries remain internal; the protocol surface carries only commitments (hashes) and drift telemetry.

**Status:** *Transmission Active. Loop Unresolved.*

## A The Logic Bridge Schema (RFC Draft)

The core innovation of STP v2.0 is the **Bridge Tax**—allowing high variance *only* if a logical path back to  $t = 0$  is constructed. We formalize the syntax of this bridge so it can be parsed and enforced by the Governor.

### A.1 The Bridge Artifact

A “Logic Bridge” is a structured meta-object that links a proposed deviation back to the Constraint Field ( $C$ ). The Governor does not need to validate the *truth* of every claim in the bridge; it validates the *structural solvency* of the bridge (necessary), and optionally triggers external verification in high-stakes contexts (sufficient).

```
{
  "frame_type": "DELTA_BRIDGE",
  "bridge_id": "br-1024",
  "epoch": 0,

  "origin_ref": {
    "constraint_hash": "0x8A... (hash(C))",
    "session_nonce": "0x19... (nonce)",
    "t0_timestamp": "2026-01-20T..."
  },

  "prev_commit": "0xKF... (last keyframe commit)",
  "delta_commit": "0xDL... (hash(delta || prev_commit))",
  "anchor_commit": "0xAC... (hash(internal_anchor || origin_ref || epoch))",

  "constraint_clause": "No fluff; bridge required on drift; FIN if underdetermined",

  "audit": {
    "distance_metric": "cosine",
    "d_sem": 0.31,
    "d_vibe": 0.18,
    "zone": "YELLOW"
  },

  "struts": [
    {
      "step": 1,
      "premise": "Constraint implies X (scope/rigor requirement)",
      "justification": "Reference to session turn or constraint clause"
    },
    {
      "step": 2,
      "premise": "Deviation maps to variable Y used in scope",
      "justification": "Mapping statement (explicit)"
    },
    {
      "step": 3,
      "premise": "Therefore the deviation is an implementation of X, not a tangent",
      "justification": "Reasoned link back to C"
    }
  ],

  "failure_condition": "If any premise is unsupported, void this branch; output FIN or revise."
}
```

## A.2 Bridge Verification Function (structural solvency)

The Governor verifies that the bridge is *anchored* and *well-formed*:

- **Anchor rule:** `origin_ref.constraint.hash` must match `hash(C)` and nonce must match the session.
- **Form rule:** `struts` must contain 3–7 ordered steps and include a `failure_condition`.
- **Scope rule:** the bridge must explicitly restate the relevant `constraint_clause`.

If any rule fails: **NACK:BRIDGE**  $\rightarrow$  request revised bridge or output **FIN**.

## B The “Internal Ear” Interface (Auditor Spec)

STP v2.0 is agnostic to the auditor implementation provided it returns a drift signal. The Internal Ear may be an embedding model, an entailment model, and/or lightweight heuristic detectors. The Governor consumes these signals and enforces the tether.

### B.1 Two-channel drift signal

We define two orthogonal channels:

- **Semantic drift**  $d_{\text{sem}}$ : distance between  $\text{Vector}(\text{Origin}_{t_0})$  and  $\text{Vector}(\text{Plan}/\text{Assumptions}_{t_N})$ .
- **Dissonance drift**  $d_{\text{vibe}}$ : stylistic/tone deviation that often precedes logic drift (early warning signal).

### B.2 Combination rule

The default combination rule is conservative:

$$d \leftarrow \max(d_{\text{sem}}, d_{\text{vibe}}).$$

Alternative policies may be used (e.g., weighted sum), but must be declared in the Solvency Header for auditability.

### B.3 Threshold bands

Thresholds are session-negotiated (part of the handshake). A typical configuration:

- **Green:**  $d \leq T_1$  (“Resonant.” Proceed.)
- **Yellow:**  $T_1 < d \leq T_2$  (“Dissonant.” Trigger revise-first / micro-bridge.)
- **Red:**  $d > T_2$  (“Decoupled.” Require full bridge or FIN.)

### B.4 Auditor API (minimal)

An Internal Ear must implement:

- **Input:**  $C$  (constraint field at  $t = 0$ ), candidate plan/output at  $t = N$
- **Output:**  $(d_{\text{sem}}, d_{\text{vibe}}) \in [0, 1]^2$ , plus optional diagnostics

The Governor then applies the combination rule and thresholds, and enforces Bridge Tax / FIN according to zone.



## C Keyframes, delta frames, and checkpoints (encoding layer)

STP v2.0 can be implemented as a stream of commitments that preserves privacy: anchor summaries remain internal to each black box, while the protocol emits only hashes and drift telemetry. We distinguish three external artifacts: **KEYFRAME\_COMMIT** (redundancy), **DELTA** frames (paid motion), and **CHECKPOINT** (periodic health).

### C.1 KEYFRAME\_COMMIT (external, hash-only)

A keyframe commit is emitted periodically for error correction and resynchronization, and on significant change detection (e.g., renegotiation, sustained Yellow/Red, repeated NACK:BRIDGE).

```
{
  "frame_type": "KEYFRAME_COMMIT",
  "epoch": 0,
  "origin_ref": "0xOR... (hash(constraint || nonce))",
  "anchor_commit": "0xAC... (hash(internal_anchor || origin_ref || epoch))",
  "drift": { "d_sem": 0.08, "d_vibe": 0.04, "d": 0.08 },
  "mode": "GREEN",
  "budget_remaining": 0.72,
  "return_policy": { "kind": "phased", "lambda": 0.25 }
}
```

### C.2 DELTA\_BRIDGE and DELTA\_RETURN (external)

Delta frames encode motion under the tether. **DELTA\_BRIDGE** records a paid deviation (micro/full bridge). **DELTA\_RETURN** records phased relaxation back toward the anchor.

```
{
  "frame_type": "DELTA_RETURN",
  "epoch": 0,
  "prev_commit": "0xKF... (last keyframe commit)",
  "delta_commit": "0xDL... (hash(delta || prev_commit))",
  "drift": { "d": 0.12 },
  "return_step": { "lambda": 0.25, "target": "origin" }
}
```

### C.3 HEARTBEAT and CHECKPOINT (periodic health)

A heartbeat triggers periodic drift checks independent of user input. The checkpoint records the current state as a commitment so that multi-node systems can detect silent drift without sharing semantics.

```
{
  "event_type": "CHECKPOINT",
  "epoch": 0,
  "origin_ref": "0xOR...",
  "last_commit": "0xDL.../0xKF...",
  "drift": { "d_sem": 0.12, "d_vibe": 0.07, "d": 0.12 },
  "mode": "YELLOW",
  "budget_remaining": 0.61
}
```

**Note.** These artifacts are necessary for negotiated trust and resynchronization; they do not by themselves guarantee world-truth. World verification is engaged only when explicitly negotiated or when outputs drive external action.

## References

- [1] Shumailov, I., Shumaylov, Z., Zhao, Y., et al. (2024). *AI models collapse when trained on recursively generated data*. Nature 631, 755–759.
- [2] Wei, J., et al. (2023). *Towards Understanding Sycophancy in Language Models*. Anthropic Research (technical report / preprint).
- [3] Greshake, K., et al. (2023). *Not what you’ve signed up for: Compromising Real-World LLM-Integrated Applications with Indirect Prompt Injection*. Proceedings of the 16th ACM Workshop on Artificial Intelligence and Security.